

OCM3D: Object-Centric Monocular 3D Object Detection

Liang Peng^{1,2} Fei Liu² Senbo Yan^{1,2} Xiaofei He^{1,2} Deng Cai¹
¹State Key Lab of CAD&CG, Zhejiang University
²FABU Inc.

{pengliang, senboyan}@zju.edu.cn {xiaofei_h, dengcai78}@qq.com {liufei}@fabu.ai

Abstract

Image-only and pseudo-LiDAR representations are commonly used for monocular 3D object detection. However, methods based on them have shortcomings of either not well capturing the spatial relationships in neighbored image pixels or being hard to handle the noisy nature of the monocular pseudo-LiDAR point cloud. To overcome these issues, in this paper we propose a novel object-centric voxel representation tailored for monocular 3D object detection. Specifically, voxels are built on each object proposal, and their sizes are adaptively determined by the 3D spatial distribution of the points, allowing the noisy point cloud to be organized effectively within a voxel grid. This representation is proved to be able to locate the object in 3D space accurately. Furthermore, prior works would like to estimate the orientation via deep features extracted from an entire image or a noisy point cloud. By contrast, we argue that the local RoI information from the object image patch alone with a proper resizing scheme is a better input as it provides complete semantic clues meanwhile excludes irrelevant interferences. Besides, we decompose the confidence mechanism in monocular 3D object detection by considering the relationship between 3D objects and the associated 2D boxes. Evaluated on KITTI, our method outperforms state-of-the-art methods by a large margin. The code will be made publicly available soon.

1. Introduction

3D object detection is of great concern in autonomous driving, and a variety of methods have been proposed. These methods can be broadly divided into two categories: LiDAR-based methods [42, 15, 31] and camera-based methods [25, 29, 19, 17, 5]. LiDAR-based methods usually offer satisfactory performances as the LiDAR point cloud can provide accurate depth measurement of the scene. However, the limited working range and the high price of LiDAR devices are the main disadvantages of this category of methods. As an alternative, camera-based methods, al-

though still far from satisfactory in terms of the detection performance, rely only on camera sensors that are much cheaper and maturer. Therefore, these approaches, especially the monocular ones [3, 24, 28, 1, 18], have drawn increasing attention from both industry and academia.

Previous image-only based monocular methods lack explicit knowledge about the depth dimension, having difficulties in predicting objects in 3D space precisely. To this end, many recent state-of-the-art monocular 3D object detection methods [8, 36, 37, 23] utilize the estimated depth map to make use of depth information. Some of them [36, 37] convert the depth map to a pseudo-LiDAR point cloud and then perform detection on it. However, most pseudo-LiDAR based methods use existing 3D detectors designed purposely for the accurate LiDAR point cloud, failing to capture the character of the highly noisy monocular pseudo-LiDAR with the long-tail defect and thus resulting in suboptimal performance.

To address the above problems, we propose a novel object-centric voxel representation. For every 2D object proposal generated by the 2D detector, we use the estimated depth map to convert pixels within the 2D box into a point cloud and build adaptive voxels on it, whose sizes depend on the 3D spatial distribution of points. With this voxel building scheme, the region with dense points that belongs more likely to an object will be partitioned with dense voxels, making more information to be encoded. On the contrary, the region with sparse points that are more likely outliers will be partitioned with sparse voxels, enabling more irrelevant interferences to be eliminated. Compared to pseudo-LiDAR, this novel voxel representation is tailor-made for monocular imagery and against the highly noisy point cloud, encoding information from the transformed point cloud and the RGB image effectively.

Another contribution is the investigation concerning the input for orientation estimation. Orientation plays a critical role in 3D object detection, tracking, and trajectory prediction. Prior works usually perform orientation estimation on deep features extracted from an entire image [29, 34, 1, 8] or a point cloud [36, 37, 23]. However, the orientation only

hinges on the object’s appearance on the local image. Semantics from the region outside the object is unnecessary, which may interfere or even overwhelm the local orientation features that matter. For this reason, we argue that the object image patch alone cropped by a 2D box is preferable as input for orientation prediction. We are the first one that points out this. Also, as detailed in Sec. 3.3, the cropped image patch needs to be properly processed to eliminate orientation ambiguity. Experiments show that our method performs best, and even comparably to the method [27] relying on LiDAR input on the pedestrian category.

Besides, an object with high 2D detection confidence may be hard to be located accurately in 3D space, *e.g.*, occluded, truncated, or distant objects, suggesting that the 2D confidence should not be naively employed in 3D detection. However, for 3D detection methods based on a 2D detector, it is hard to learn 3D confidence without explicit labels during training. Consequently, most SOTA methods [36, 37, 23] directly apply the 2D confidence as the 3D detection confidence score. To resolve this problem, we propose to decompose the 3D confidence mechanism into the confidence in 2D detection and the lifting hardness of an object from 2D to 3D, where the former can be easily obtained from a 2D detector. The lifting hardness is measured by the relationship between the 2D box and 3D box projections as we believe that objects are worthy of high confidences when the 2D and 3D detections fit tightly on the image plane. Note that the proposed method can be plugged into any 2D-detector-based method without training and significantly boost the performance.

Finally, we discover that the training set in Eigen split [9] that is widely adopted for training depth estimation networks overlaps the KITTI 3D object detection validation set. This data leakage results in that almost all the depth-based monocular 3D object detection methods are overrated in the KITTI validation set. To remedy this problem, we introduce a new dataset split. In conclusion, our main contributions can be summarized as follows:

- We propose a novel object-centric voxel representation, which effectively encodes the noisy point cloud and the RGB image.
- We point out that a cropped object image patch is preferable as input compared to an entire image for orientation prediction.
- A novel decomposed 3D confidence mechanism is designed for monocular 3D object detection, taking both 2D confidence and the lifting hardness from 2D to 3D into consideration.
- A widely existing data leakage issue is pointed out and remedied. Also, extensive experiments show that we

achieve state-of-the-art performance on KITTI monocular 3D detection benchmark with a significant margin.

2. Related Work

LiDAR-based Methods: Most state-of-the-art 3D object detection methods [32, 15, 11, 40, 33, 38, 39, 12] employ LiDAR as it can provide accurate point clouds, in which voxel-based methods are utmost related to our method. VoxelNet [42] divides the LiDAR point cloud into a voxel grid with a fixed voxel size, a group of points that belong to the same voxel is fed into a fully connected network to form the unified feature representation. The 2D convolution of these higher-order voxel features is carried out to obtain detection results. Following voxel-based approaches [35, 20] follow this line of thought. Such designs tailored for accurate LiDAR points are not suitable for monocular methods since the point cloud transformed from the predicted depth map is much more inaccurate. Points of the same object distribute chaotically within the whole 3D space (*e.g.*, 3D points transformed from neighboring pixels can appear far away in 3D space, even they are supposed to be much close). CNNs are hard to capture local clues from voxels built with the way above.

Representations for Monocular Methods: Monocular methods can be roughly divided into image-only based methods and depth map based methods according to representations. Previous image-only based methods usually take prior knowledge as guidance in training or as constraints in post-processing. For example, Mono3D [3] uses semantic and context priors as guidance to generate 3D proposals, require extra semantic and instance segmentation networks. M3D-RPN [1] uses different convolution kernels in row-spaces, trying to explore different features in different depth ranges and RTM3D [18] predicts perspective key points to refine initial guess by solving a constrained optimization problem. These image-only based methods have difficulties in exploiting the hidden depth information, remaining room for improvement. To make use of depth information explicitly, many methods use the depth map produced by a depth estimating network. D4LCN [8] uses different kernels generated by depth maps, but not full use of the explicitly spatial relationship. Pseudo-LiDAR [36, 37] converts image-based depth maps to point clouds to mimic the LiDAR signal, and then directly using LiDAR-based 3D detectors. On this basis, AM3D [23] augments this method by embedding RGB values to generate attention maps for multi-modal features fusion. However, pseudo-LiDAR based methods ignore the significant gap between the transformed noisy point cloud and the accurate LiDAR point cloud. Furthermore, PatchNet [22] points out that the effectiveness of pseudo-LiDAR comes from the coordinates

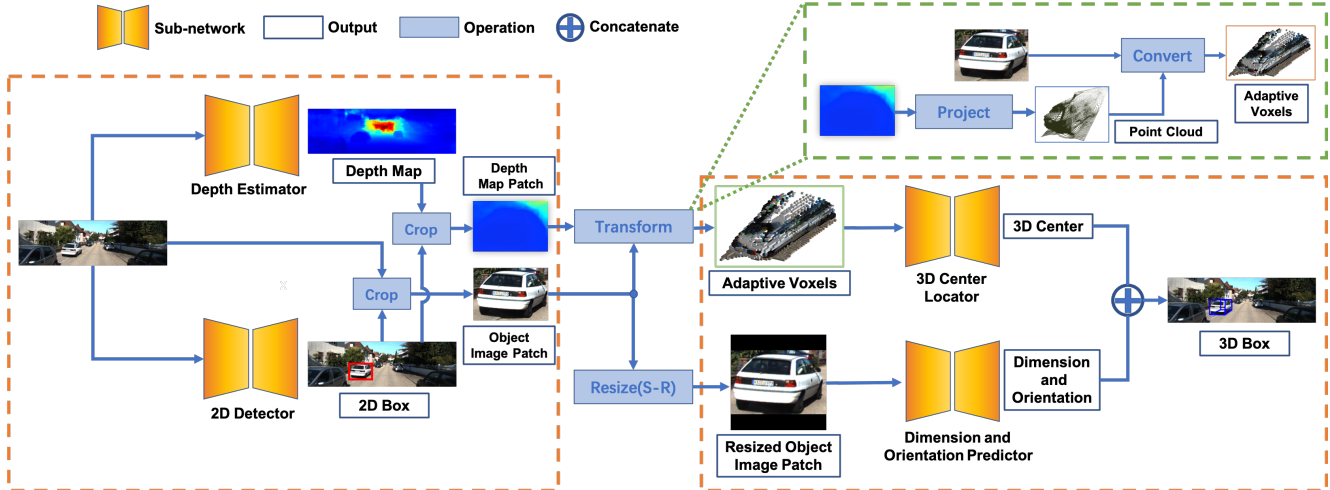


Figure 1: OCM3D framework. OCM3D consists of four subnetworks for 2D detection, depth estimation, 3D center location, and dimension and orientation prediction, respectively. The 2D boxes and depth map are firstly produced by off-the-shelf networks. For each object, the estimated depth map and 2D box are used to construct adaptive voxels (Sec. 3.2), on which a 3D center locator predicts the object’s 3D center. By utilizing another subnetwork, the object’s dimension and orientation are predicted from the cropped image patch (Sec. 3.3). Best viewed in color with zoom in.

transform. It organizes the pseudo-LiDAR as the image representation and utilizes existing 2D CNN designs to extract deep features. However, it does not consider utilizing the 3D spatial relationship inside the object.

Orientation Estimation: Many approaches use the entire image directly for orientation estimation, such as [29, 34, 1, 8]. The global semantics extracted from the region outside the object may interfere or even overwhelm the local semantic that matters for orientation estimation. Pseudo-LiDAR [36, 37] performs prediction on the point cloud, which is hard to infer orientations of objects with noisy points. Deep3DBBox [25] takes cropped object image patches as input, using MultiBin architecture for orientation regression. However, they process the image patch by using naive resizing operation, making the orientation ambiguous. We also utilize the object image patch and introduce another simple yet effective resizing scheme to replace the naive resizing without damaging any semantic clues. Based on this proper resizing scheme, we argue that the object image patch alone provides adequate features to predict orientations.

Monocular 3D Object Detection Confidence: Previous works in monocular 3D object detection realize the need for 3D box confidence. FQNet [19] employs 3D IoU loss calculated from the predicted 3D box and the ground truth, mainly considering the 3D object dimensions. MonoDIS [34] takes advantage of the 3D box loss, represents the confidence of 3D detection. As a video-based method, the recent concurrent work, Kinematic 3D [2], introduces a

self-balancing confidence loss, re-balancing hard 3D boxes and focusing on relatively achievable samples. All of these methods consider the confidence towards loss function. By contrast, our method takes the physical projections into consideration, gains promising performances, and can be adopted conveniently by other methods.

3. Methods

3.1. Overview

Our object-centric monocular 3D detection framework (OCM3D) can be divided into two stages. As shown in Fig. 1, stage one contains 2D detection and depth map estimation, producing 2D boxes and the image-based depth map by two off-the-shelf networks [27, 16], as a common practice [36, 37, 23, 22]. Our method mainly focuses on stage two that comprises object 3D center location and dimension & orientation prediction.

In stage two, dimensions and orientations are predicted via cropped object image patches with a proper resizing scheme (Sec. 3.3). At the same time, with the addition of the depth map, adaptive voxels are constructed (Sec. 3.2) and then fed into a 3D fully convolutional network [21, 30, 7] (see the supplementary material). The network outputs a 3D heat map to obtain the 3D center, since we formulate the 3D center location problem as a key-point prediction task (Sec. 3.2). Finally, the 3D detection confidence is determined by the 2D detection confidence and the lifting hardness from 2D to 3D (Sec. 3.4).

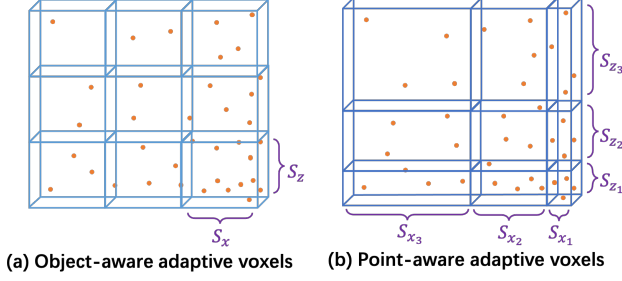


Figure 2: Illustration of different adaptive voxels. Only X and Z axis are shown for convenience. Orange dots represent the ROI point cloud, S_x, S_z refer to the voxel sizes along X and Z axis, and each small cube is a voxel.



Figure 3: Illustration of obtaining the 3D heatmap.

3.2. Adaptive Voxels

In this section, we describe how to construct adaptive voxels. Unlike prior works that divide the entire 3D space into a voxel grid with the fixed voxel size, we first partition the 3D space to obtain an ROI 3D space, which contains only the specified object, enabling fine-grained voxels to be constructed for encoding detailed features. Thanks to the fine-grained design, RGB values on the image can be utilized directly as initial features for corresponding voxels via the projection, allowing the network to utilize the object’s visual clues and spatial relationship adequately. Therefore the remaining problem is how to customize the voxel size for every object to fit their noisy and varied point cloud distribution.

Specifically, as shown in Fig. 1 and Fig. 2, given a 2D object proposal in an image, the ROI point cloud is derived through the estimated depth map. Then the cropped ROI 3D space can be divided into a voxel grid of shape $n_x \times n_y \times n_z$, along the X, Y, Z axis, respectively, where voxel sizes are determined according to:

$$S_{i \in \{x, y, z\}} = \frac{\max(P_i) - \min(P_i)}{n_i}, \quad (1)$$

in which i refers to the axis, S is the voxel size, and P is the 3D coordinates of the ROI point cloud. We can see that $(\max(P_i) - \min(P_i))$ represents 3D spatial offsets along each axis in the ROI point cloud. It indicates that voxel sizes are determined by the local spatial distribution and can

adapt online to different objects to encode information better. In light of this, we term it as object-aware adaptive size (O-A).

Nevertheless, the object-aware adaptive size is hard to explore the in-depth spatial distribution inside objects. In the noisy ROI point cloud, points can be dense somewhere while sparse elsewhere, resulting in an unbalanced feature distribution in the voxel grid. To effectively handle this, based on object-aware adaptive size, another type of adaptive voxel size is introduced, where the size of each voxel in the voxel grid can be different to fit the data distribution.

$$\begin{cases} MX_{k \in \{0, 1, \dots, n_x\}} = \text{sort}(P_x) \left[\left\lfloor k \frac{N}{n_x} \right\rfloor \right] \\ MY_{k \in \{0, 1, \dots, n_y\}} = \text{sort}(P_y) \left[\left\lfloor k \frac{N}{n_y} \right\rfloor \right] \\ MZ_{k \in \{0, 1, \dots, n_z\}} = \text{sort}(P_z) \left[\left\lfloor k \frac{N}{n_z} \right\rfloor \right] \end{cases} \quad (2)$$

As shown in Eq. 2 and Fig. 2, sort is the sort operation among 3D coordinates (e.g., $\text{sort}(P_x)$ refers to the sorted x coordinates of the ROI point cloud), N is the number of points, and $\lfloor \cdot \rfloor$ denotes gathering the corresponding index. MX represents the x coordinates with given indexes in the sorted ROI point cloud. MY and MZ follow this line. More specifically, MX, MY, MZ refer to retrieve a series of coordinates that locate on specified indexes, dividing sorted coordinates into n_x, n_y, n_z parts. Therefore, the size of each voxel in the voxel grid depends on the spatial offset between head elements of two adjacent sorted parts (e.g., the k th voxel size along X axis is $MX_{k+1} - MX_k$), making it varies a lot inside the unevenly distributed ROI point cloud. The region with dense points more likely belongs to an object, and it will be partitioned with small voxel sizes, enabling more information to be maintained. By contrast, the region with sparse points that are more likely outliers will be partitioned with large voxel sizes, generating sparse voxels, allowing more irrelevant interferences to be eliminated. We call this type of voxel size the point-aware adaptive size (P-A).

Aiming to obtain the 3D object center, as shown in Fig. 3, the location network acquires the adaptive voxels and outputs the location heatmap with the same shape of input, where the 3D object center is the location with max probability. The heatmap groundtruth is generated as described in [41], and smooth L1 loss is adopted during the training. Also, to avoid the object-aware adaptive size is heavily affected by outlier points, we remove these points whose depths are much larger than the average depth [23].

3.3. Input for Orientation Estimation

The orientation in 3D object detection contains the global orientation and the observation angle (local orientation), where the former is defined under the world space and

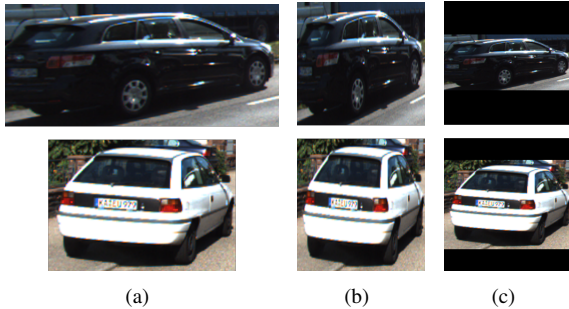


Figure 4: Illustration of different resizing schemes. (a): the original cropped object image patch (OIP), (b): OIP after naive resizing, (c): OIP after Shape-Retaining(S-R).

the observation angle is defined under the camera coordinates. The global orientation is the goal and can be derived from the observation angle and the camera viewing angle.

Most state-of-the-art methods [1, 28, 23, 8] directly infer the observation angle via deep features extracted from an entire image or a noisy point cloud. However, the observation angle can be better predicted via the local RoI information as it only hinges on the object image patch. Information from context does not offer much help because semantics from the region outside the object may bring irrelevant or even misleading information for local semantics. Also, the noisy point cloud can not express the object orientation precisely due to the noisy nature. Thus we argue that the cropped image patch alone is preferable as the input for orientation estimation.

Interestingly, some prior works such as [25] also utilize the cropped image patch to achieve the orientation. However, they are unaware of the underlying mechanism and improperly use a naive resizing scheme to resize object image patches to fit the network input. As shown in Fig. 4, naive resizing may destroy the visual semantics and thus tends to make the observation angles ambiguous. To this aim, we introduce another simple yet angle-preserving resizing scheme, called Shape-Retaining (S-R), i.e., resizing to match larger edge while keeping aspect ratio and then padding zeros symmetrically to get the desired shape. In this way, information for the observation angle is maintained without damaging, thus better prediction accuracy can be achieved, as shown in Tab. 5 and Tab. 7. We obviously outperforms other monocular methods and even gives comparable results to the LiDAR-based method [27] on the pedestrian category.

3.4. Decomposed 3D Detection Confidence

In 2D detection, the confidence usually denotes the classification score to rate a box, which can be naturally learned

in 2D detection. An object assigned with a high confidence score means that it can be easily discriminated from the 2D image, but does not necessarily mean that it can be easily located accurately in 3D space. For example, truncated/occluded cars, which can be easily detected in 2D, are considered as hard cases in 3D. Therefore, directly adopting 2D confidence as 3D confidence, like many prior works [36, 37, 23], are not suggested. Also, the 3D detection confidence can be learned during training by designing a specific loss using the 3D IoU [19], but this manner is indirect and ignores the geometry relationship of 3D boxes and associated 2D boxes. To make estimation of 3D detection confidence more accurately, we decompose it as follows:

$$Conf_{3D} = Conf_{2D} \times Conf_{2D \rightarrow 3D} \quad (3)$$

The former is just the 2D detection confidence, and the latter measures the lifting hardness of the object from 2D to 3D, which considers how difficulty in locating a known 2D object in 3D space. Specifically, we achieve the lifting hardness by using the 2D-3D box IoU as the fact that the precise 3D box projections should fit its corresponding 2D box tightly, which is a much simpler yet more effective method. In particular, with the growth of the object distance, we find that the 2D-3D box IoU becomes larger for the object with the same location error, which means that the distance should also be taken into consideration. Fig. 5 shows some examples. Based on the above analysis, the proposed 3D detection confidence is calculated as follows:

$$Conf_{2D \rightarrow 3D} = \frac{IoU(box_{3D-proj}, box_{2D})}{e^{dis/\lambda}} \quad (4)$$

in which dis refers to the distance from the predicted object to the camera optical center, i.e., $\sqrt{x^2 + y^2 + z^2}$, $Conf_{2D \rightarrow 3D}$ denotes the lifting hardness from 2D to 3D, λ is used to scale the depth, and $box_{3D-proj}$ refers to the 2D bounding box on predicted 3D box projections. When the distance becomes larger, the estimation error of location is increasing exponentially (see the supplementary material) as the information for location is less rich. Therefore, we adopt the exponential transformation for the distance in Eq. 4 to account for location error implicitly. In this way, we balance influences brought from the 2D-3D box IoU and the location. It is worthy to note that our method can be plugged into any other 2D-detector-based monocular method without cost to boost the performance.

4. Experiments

4.1. Implementation Details

We implement our framework using Pytorch [26] and on a Nvidia 1080Ti GPU, applying Adam optimizer [13], in which the learning rate is set to 10^{-4} . The encoder-decoder network architecture is employed for the location



Figure 5: IoU of 3D boxes on image plane at different distances. Top left refers to the BEV map and the boxes in the RGB image are the projections. With the growth of distance, the IoU of 3D boxes on the image plane becomes larger, although their offset keeps unchanged (2 meters here).

Approach	$AP_{BEV}/AP_{3D} (IoU=0.5) _{R_{11}}$			$AP_{BEV}/AP_{3D} (IoU=0.7) _{R_{11}}$		
	Easy	Moderate	Hard	Easy	Moderate	Hard
OFTNet [29]	-	-	-	11.06/4.07	8.79/3.27	8.91/3.29
RoI-10D [24]	-	-	-	14.50/10.25	9.91/6.39	8.73/6.18
MonoDIS [34]	-	-	-	24.26/18.05	18.43/14.98	16.95/13.42
MonoGRNet [28]	-/50.51	-/36.97	-/30.82	-/13.88	-/10.19	-/7.62
Deep3DBBox [25]	30.02/27.04	23.77/20.55	18.83/15.88	9.99/5.85	7.71/4.10	5.30/3.84
Mono3D [3]	30.50/25.19	22.39/18.20	19.16/15.52	5.22/2.53	5.19/2.31	4.13/2.31
FQNet [19]	32.57/28.16	24.60/21.02	21.25/19.91	9.50/5.98	8.02/5.50	7.71/4.75
MonoPSR [14]	56.97/49.65	43.39/41.71	36.00/29.95	20.63/12.75	18.67/11.48	14.45/8.59
M3D-RPN [1]	55.37/48.96	42.49/39.57	35.29/33.01	25.94/20.27	21.18/17.06	17.90/15.21
Pseudo-LiDAR [36]	<u>63.42/56.54</u>	40.86/37.53	37.69/32.24	31.88/ <u>24.12</u>	20.84/15.74	18.92/14.96
D4LCN [8]	54.35/51.30	40.33/35.10	33.96/32.46	26.00/19.38	20.73/16.00	17.46/12.94
RTM3D [18]	57.47/54.36	<u>44.16/41.90</u>	<u>42.31/35.84</u>	25.56/20.77	<u>22.12/16.86</u>	<u>20.91/16.63</u>
PatchNet [22]	64.87/58.63	40.94/38.20	38.16/32.50	<u>32.30/25.76</u>	21.25/ <u>17.72</u>	19.04/15.62
Ours	<u>61.78/58.07</u>	45.65/42.61	42.32/36.62	33.24/23.65	24.25/17.75	21.67/15.93

Table 1: Comparison of our method to other monocular 3D object detection methods for car category on KITTI validation set. For a fair comparison, depth map based methods [36, 8, 22] are retrained using our proposed depth dataset split. To include more methods in comparison, all the methods are evaluated with metric $AP|_{R_{11}}$.

sub-network and is trained with a total of 50 epochs (see the supplementary material for more details). The location results are refined by the 2D-3D box consistency. For the orientation and dimension sub-network, we use the protocol introduced in [25]. Depth estimator and 2D detector are brought from [16] and [27], using our remedied depth dataset split to train the depth estimator. The transformed point cloud is rotated as described in [27], and the shape of the voxel grid (n_x, n_y, n_z) is set to (32, 16, 64). λ in Eq. 4 is 80 in default since 80 meters is the usual max depth in KITTI. We provide extra experimental results in the supple-

mentary material due to the space limitation.

4.2. Dataset and Metrics

KITTI Dataset: The KITTI dataset is widely used for monocular 3D object detection. KITTI [10] provides 7,481 images for training and 7,518 images for testing. The test set is confidential and can only be tested on the KITTI website, while the training set is publicly available. To keep consistent with previous works, we adopt the dataset split described in [4] for comprehensive experiments. The available data is divided into a training set (3,712 images) and a

O-A	S-R	3D-Conf	P-A	AP_{BEV}/AP_{3D} (IoU=0.7) $ _{R_{40}}$		
				Easy	Moderate	Hard
✓				22.12/12.37	11.89/6.08	8.77/4.41
✓	✓			24.45/14.60	12.78/7.20	9.37/5.31
✓	✓	✓		27.93/17.36	16.52/10.04	12.04/7.25
✓	✓	✓	✓	28.13/17.57	19.90/ 12.10	16.86/ 10.36
✓	✓	✓	✓	29.57/18.08	20.01/12.09	16.89/10.15

Table 2: Ablation study.

3D Confidence	AP_{BEV}/AP_{3D} (IoU=0.7) $ _{R_{40}}$		
	Easy	Moderate	Hard
Learning-based	26.85/15.39	16.41/9.33	14.70/9.02
Projection-based	29.57/18.08	20.01/12.09	16.89/10.15

Table 3: Comparisons on different 3D confidence schemes.

validation set (3,769 images).

Depth Dataset Split for 3D Object Detection: Existing depth estimation methods usually are trained upon the split proposed in [9], in which 32 scenes among all KITTI scenarios are used for training. To the best of our knowledge, most depth-based monocular 3D object detection methods estimate depth directly from the existing pre-trained models. However, the training set for depth estimation overlaps with the validation set for 3D object detection, resulting in overrating these methods. To eliminate this problem, we introduce a new depth dataset split for 3D object detection. Specifically, we exclude scenes that emerge in the KITTI 3D object validation set. The remaining scenes consist of the new training set, and others are divided into the validation set. This depth dataset split contains data from non-overlapping sequences and fixes the data leakage issue. Thus we re-train and re-evaluate depth-based methods [36, 8, 22] using the official publicly available code.

Evaluation Metrics: We conduct experiments on the KITTI validation set and official test set under two core tasks: bird’s eye view (BEV) and 3D object detection in three difficulties. Difficulties of objects are subdivided into easy, moderate, and hard according to the occlusion level, truncation, and bounding box height. Many previous works evaluate their results with AP_{11} , while AP_{40} proposed in [34] is suggested to be adopted by the official KITTI benchmark recently. Therefore we provide performances under both AP_{11} and AP_{40} metrics on the car category to comprehensively evaluate the proposed method.

4.3. Quantitative Results

We compare our method with current state-of-the-art monocular methods in the KITTI validation set. Note that depth map based methods are trained upon the new depth

Approach	Easy	Moderate	Hard
AM3D [23]*	25.03/16.50	17.32/10.74	14.91/9.52
D4LCN [8]*	22.51/16.65	16.02/11.72	12.55/9.51
PatchNet [22]*	22.97/15.68	16.86/11.12	14.97/10.17
ROI-10D [24]	9.78/4.32	4.91/2.02	3.74/1.46
MonoGRNet [28]	18.19/5.74	11.17/9.61	8.73/4.25
MonoPSR [14]	18.33/10.76	12.58/7.25	9.91/5.85
M3D-RPN [1]	21.02/14.76	13.67/9.71	10.23/7.42
MonoPair [6]	19.28/13.04	14.83/9.99	12.89/8.65
RTM3D [18]	19.17/14.41	14.20/10.34	11.99/ 8.77
Ours	27.87/17.48	17.13/10.44	13.53/7.87

Table 4: Comparisons on KITTI testing set. Note that * denotes the method use a different depth estimator with different depth dataset split compared to ours.

dataset split. As shown in Tab. 1, our method achieves new state-of-the-art results. For the BEV and 3D tasks, we gain significant improvements compared to other state-of-the-art results. Also, we provide comparisons on KITTI testing set as shown in Tab. 4, in which our method still achieves state-of-the-art. It is worthy to note that results of our method should not be compared directly with other depth map based methods which utilize a different depth estimator trained with the un-remedied depth dataset split.

4.4. Detailed Analysis

Ablation Study: We provide extensive experiments to study the impact of each critical component in our framework. As shown in Tab. 2, modules: object-aware adaptive voxel size (O-A), point-aware adaptive voxel size (P-A), Shape-Retaining (S-R), decomposed 3D confidence (3D-Conf) are gradually added to the framework, producing growing improvements. We also compare different decomposed 3D confidence schemes mentioned in Sec. 3.4, i.e., the learning-based method using 3D IoU [19] and our 2D-3D projection-based method. We can see that the simple yet effective projection-based method performs much better.

Extensibility of Decomposed 3D Confidence Scheme:

As mentioned in Sec. 3.4, our decomposed 3D confidence scheme only requires the 2D detection confidence and the predicted 3D detection results, meaning that it can be easily plugged into any monocular 3D object detection method based on a 2D detector. As shown in Tab. 6, by employing our decomposed 3D confidence scheme, Pseudo-LiDAR [36] gains **39.87%/47.4%** relative improvements, and PatchNet [22] gains **33.63%/40.76%** relative improvements in *moderate* setting. This significant boosting demonstrates the effectiveness of our decomposed 3D confidence scheme.

Different Inputs for Orientation Estimation: In Sec. 3.3,

Approaches/Input	Dataset	AOS $_{R40}$			AP $_{BEV}$ $_{R40}$			AP $_{3D}$ $_{R40}$		
		Easy	Mod	Hard	Easy	Mod	Hard	Easy	Mod	Hard
M3D-RPN [1]/Entire image	Validation set	44.86	37.57	31.72	5.18	4.48	3.61	4.75	3.55	2.79
M3D-RPN/Image patch(Ours)		63.01(+18.15)	53.11(+16.54)	44.69(+12.97)	6.02(+0.84)	4.79(+0.31)	4.03(+0.42)	5.24(+0.49)	4.33(+0.78)	3.63(+0.84)
M3D-RPN [1]/Entire image	Test set	44.33	31.88	28.55	5.65	4.05	3.29	4.92	3.48	2.94
M3D-RPN/Image patch(Ours)		51.90(+7.57)	37.78(+5.90)	33.95(+5.40)	6.53(+0.98)	4.46(+0.39)	4.10(+0.91)	5.70(+0.78)	4.11(+0.63)	3.37(+0.43)
F-PointNet [27]/LiDAR	Validation set	81.18(+1.05)	72.06(+0.01)	66.00	72.15(+1.79)	62.16(+1.54)	54.13(+0.80)	64.99(+1.43)	55.12(+1.65)	47.26(+0.87)
F-PointNet/Image patch(Ours)		80.13	72.05	66.05(+0.05)	70.36	60.52	53.33	63.56	53.47	46.39

Table 5: Comparison of different inputs for estimating orientation on the pedestrian category on KITTI benchmark. Our method gives much better results and even be comparable with the LiDAR-based method [27].

Methods	AP $_{BEV}$ /AP $_{3D}$ (IoU=0.7) $_{R40}$		
	Easy	Moderate	Hard
P-LiDAR [36]	27.29/18.43	15.10/9.74	12.61/7.81
P-LiDAR+3D-Conf	31.77/21.48	21.12/14.36	17.73/11.67
PatchNet [22]	27.97/20.75	15.67/11.26	12.92/9.07
PatchNet+3D-Conf	31.52/23.57	20.94/15.85	17.40/12.78

Table 6: Extensibility of our decomposed 3D confidence scheme. "P-LiDAR" in the table denotes Pseudo-LiDAR [36]. The proposed 3D confidence scheme brings significant performance improvements.

Method	Metric	Easy	Moderate	Hard
M3D-RPN [1]	AOS(Org.)	89.37	81.58	65.74
	AOS(Ours)	89.42	82.53	66.64
	2D Det.	90.02	83.14	67.37
D4LCN [8]	AOS(Org.)	91.75	82.97	66.45
	AOS(Ours)	92.12	83.27	66.81
	2D Det.	92.76	84.40	67.86
Pseudo-LiDAR [36]	AOS(Org.)	92.85	82.29	78.51
	AOS(Ours)	96.15	89.84	85.56
	2D Det.	96.48	90.30	87.62

Table 7: Comparison of different way to estimate orientation on car category on the KITTI validation set. "2D Det." in the table denotes the AP of 2D detection, which is the up-bound for AOS.

we argue that the object image patch alone provide enough information for orientations of objects in the image, which is proved here via quantitative experiments. Prior state-of-the-art (SOTA) methods, M3D-RPN [1], D4LCN [8], Pseudo-LiDAR [36] are chosen as the baseline. They predict orientations via extracted features from the entire image or the point cloud. By contrast, we only use the same 2D boxes produced by them and apply our method. We show the results of Average Orientation Similarity (AOS) on the car category in Tab. 7, where AOS(Org.) denotes performances produced by the original method and 2D Det. denotes the AP of 2D detection, which is the up-bound for

AOS. Our method performs much better than the original method with a significant gap, demonstrating our perspective. Besides, we conduct experiments on the pedestrian category as shown in Tab. 5. By replacing only our orientations to M3D-RPN [1], dramatic improvements are obtained both in the validation set and test set. Even if compared to the LiDAR-based method [27], it achieves comparable performances. These huge gains can be attributed to that our method maintains the information as much as possible and makes the model focus on the object image patch. Especially for the pedestrian category, which has a high aspect ratio, previous methods bring in interferences of context, thus perform much worse.

4.5. Limitations

Our current framework still has some limitations. First, compared to other end-to-end algorithms such as [1], our training process is of complication. All sub-modules can be integrated and trained with an end-to-end learning approach. Second, we have not taken some more useful information into consideration, such as semantics. Such aspects will be explored in our future work.

5. Conclusions

In this paper, we propose a novel data representation to encode the depth map and the RGB image in which voxel sizes are adaptively determined by the spatial distribution of the transformed point cloud. This representation is used to predict objects' locations. Furthermore, we propose that the object image patch alone with a proper resizing scheme is a better input than other methods that utilize the entire image or the noisy point cloud. Experiments demonstrate our perspective. Besides, we introduce a novel decomposed 3D confidence approach, constructing the connection between 2D detection and 3D detection. Finally, we reveal a data-leakage problem and propose a new dataset split to fix it. Extensive experiments demonstrate that our method is promising and outperforms other state-of-the-art methods.

References

- [1] Garrick Brazil and Xiaoming Liu. M3d-rpn: Monocular 3d region proposal network for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9287–9296, 2019. 1, 2, 3, 5, 6, 7, 8
- [2] Garrick Brazil, Gerard Pons-Moll, Xiaoming Liu, and Bernt Schiele. Kinematic 3d object detection in monocular video. *arXiv preprint arXiv:2007.09548*, 2020. 3
- [3] Xiaozhi Chen, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. Monocular 3d object detection for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2147–2156, 2016. 1, 2, 6
- [4] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals using stereo imagery for accurate object class detection. *IEEE transactions on pattern analysis and machine intelligence*, 40(5):1259–1272, 2017. 6
- [5] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Dsgn: Deep stereo geometry network for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12536–12545, 2020. 1
- [6] Yongjian Chen, Lei Tai, Kai Sun, and Mingyang Li. Monopair: Monocular 3d object detection using pairwise spatial relationships. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12093–12102, 2020. 7
- [7] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *International conference on medical image computing and computer-assisted intervention*, pages 424–432. Springer, 2016. 3
- [8] Mingyu Ding, Yuqi Huo, Hongwei Yi, Zhe Wang, Jianping Shi, Zhiwu Lu, and Ping Luo. Learning depth-guided convolutions for monocular 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11672–11681, 2020. 1, 2, 3, 5, 6, 7, 8
- [9] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. 2, 7
- [10] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012. 6
- [11] Chenhang He, Hui Zeng, Jianqiang Huang, Xian-Sheng Hua, and Lei Zhang. Structure aware single-stage 3d object detection from point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11873–11882, 2020. 2
- [12] Tengeng Huang, Zhe Liu, Xiwu Chen, and Xiang Bai. Epnet: Enhancing point features with image semantics for 3d object detection. *arXiv preprint arXiv:2007.08856*, 2020. 2
- [13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [14] Jason Ku, Alex D Pon, and Steven L Waslander. Monocular 3d object detection leveraging accurate proposals and shape reconstruction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11867–11876, 2019. 6, 7
- [15] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12697–12705, 2019. 1, 2
- [16] Jin Han Lee, Myung-Kyu Han, Dong Wook Ko, and Il Hong Suh. From big to small: Multi-scale local planar guidance for monocular depth estimation. *arXiv preprint arXiv:1907.10326*, 2019. 3, 6
- [17] Peiliang Li, Xiaozhi Chen, and Shaojie Shen. Stereo r-cnn based 3d object detection for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7644–7652, 2019. 1
- [18] Peixuan Li, Huaici Zhao, Pengfei Liu, and Feidao Cao. Rtm3d: Real-time monocular 3d detection from object keypoints for autonomous driving. *arXiv preprint arXiv:2001.03343*, 2020. 1, 2, 6, 7
- [19] Lijie Liu, Jiwen Lu, Chunjing Xu, Qi Tian, and Jie Zhou. Deep fitting degree scoring network for monocular 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1057–1066, 2019. 1, 3, 5, 6, 7
- [20] Zhijian Liu, Haotian Tang, Yujun Lin, and Song Han. Pointvoxel cnn for efficient 3d deep learning. In *Advances in Neural Information Processing Systems*, pages 965–975, 2019. 2
- [21] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 3
- [22] Xinzhu Ma, Shinan Liu, Zhiyi Xia, Hongwen Zhang, Xingyu Zeng, and Wanli Ouyang. Rethinking pseudo-lidar representation. *arXiv preprint arXiv:2008.04582*, 2020. 2, 3, 6, 7, 8
- [23] Xinzhu Ma, Zhihui Wang, Haojie Li, Pengbo Zhang, Wanli Ouyang, and Xin Fan. Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6851–6860, 2019. 1, 2, 3, 4, 5, 7
- [24] Fabian Manhardt, Wadim Kehl, and Adrien Gaidon. Roi-10d: Monocular lifting of 2d detection to 6d pose and metric shape. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2069–2078, 2019. 1, 6, 7
- [25] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7074–7082, 2017. 1, 3, 5, 6
- [26] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming

- Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*, pages 8026–8037, 2019. [5](#)
- [27] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 918–927, 2018. [2](#), [3](#), [5](#), [6](#), [8](#)
- [28] Zengyi Qin, Jinglu Wang, and Yan Lu. Monogrnet: A geometric reasoning network for monocular 3d object localization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8851–8858, 2019. [1](#), [5](#), [6](#), [7](#)
- [29] Thomas Roddick, Alex Kendall, and Roberto Cipolla. Orthographic feature transform for monocular 3d object detection. *arXiv preprint arXiv:1811.08188*, 2018. [1](#), [3](#), [6](#)
- [30] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [3](#)
- [31] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020. [1](#)
- [32] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *arXiv preprint arXiv:1907.03670*, 2019. [2](#)
- [33] Weijing Shi and Raj Rajkumar. Point-gnn: Graph neural network for 3d object detection in a point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1711–1719, 2020. [2](#)
- [34] Andrea Simonelli, Samuel Rota Buló, Lorenzo Porzi, Manuel López-Antequera, and Peter Kotschieder. Disentangling monocular 3d object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1991–1999, 2019. [1](#), [3](#), [6](#), [7](#)
- [35] Bei Wang, Jianping An, and Jiayan Cao. Voxel-fpn: multi-scale voxel feature aggregation in 3d object detection from point clouds. *arXiv preprint arXiv:1907.05286*, 2019. [2](#)
- [36] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8445–8453, 2019. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [8](#)
- [37] Xinshuo Weng and Kris Kitani. Monocular 3d object detection with pseudo-lidar point cloud. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. [1](#), [2](#), [3](#), [5](#)
- [38] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3dssd: Point-based 3d single stage object detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11040–11048, 2020. [2](#)
- [39] Maosheng Ye, Shuangjie Xu, and Tongyi Cao. Hvnet: Hybrid voxel network for lidar based 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1631–1640, 2020. [2](#)
- [40] Jin Hyeok Yoo, Yeocheol Kim, Ji Song Kim, and Jun Won Choi. 3d-cvf: Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection. *arXiv preprint arXiv:2004.12636*, 2020. [2](#)
- [41] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. [4](#)
- [42] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4490–4499, 2018. [1](#), [2](#)